

# Lecture # 26

crash consistency  
problem & solutions.

⊛ recap:

VFS — virtual file system (Linux lingo)

- standard file model

- inode

- superblock

- dentry

- file object

in-memory objects

for file system

interface and usage in OS.

- f-ops

(file operations)  
points to

- each FS has to provide/adhere to

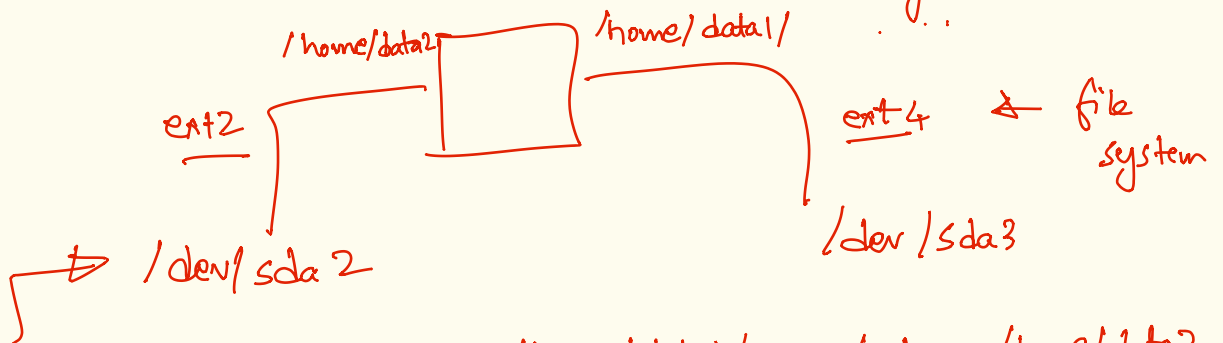
file model for OS to use/operate the VFS

- individual file systems will have their own

on-disk objects and corresponding operations

& logic.

read  
write  
open  
close  
seek



file name

which  
points to

a device

⇒ cp /home/data/peace.txt /home/data2

- copy across file systems!

mount ?

/home



/dev/sda1

ext2 ← FS

mount /dev/sda1

/home/data ← files

ext2

/home/data/

↳ maps the root

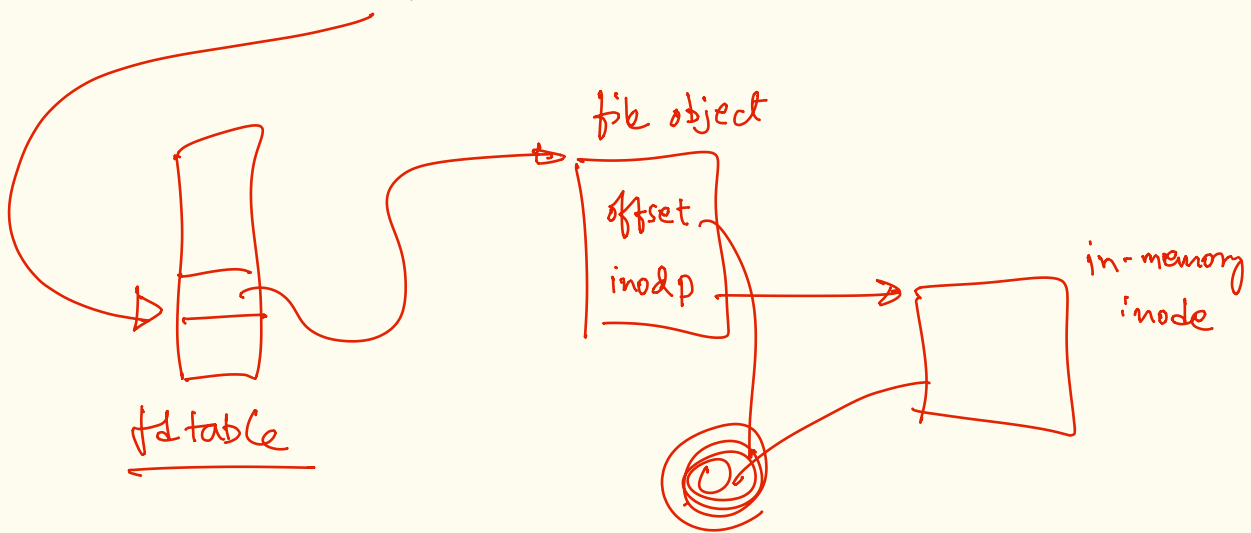
inode

on /dev/sda1

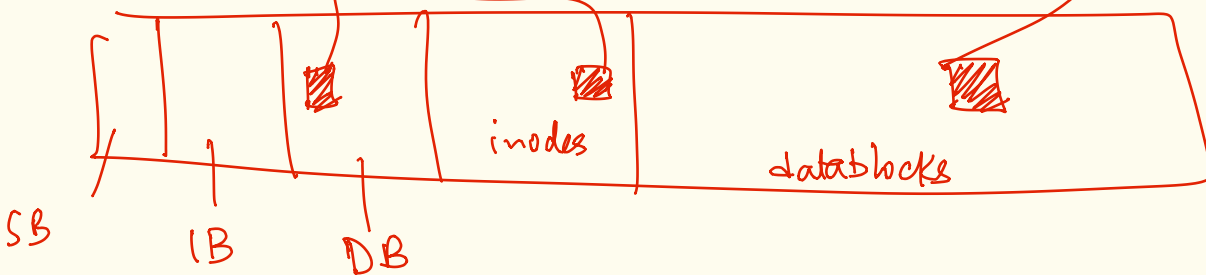
ls /home/data/

# # crash consistency problem.

```
write (fd, buf, N);
```



- new datablock reserved
- write to the datablock
- write/update the data bitmap
- write on the on-disk inode



(\*) crash scen

# # crash scenarios

write ~~to~~ to a file  $\Rightarrow$  3 writes to disk

C: consistent		<u>inode</u>	<u>data bitmap</u>	<u>data block</u>	
IC: inconsistent	C	✓	✓	✓	
	C	X	X	X	
<u>midly IC</u>	<u>data lost</u>	X	X	✓	
IC	<u>inode has bad info.</u>	✓	X	X	
IC	<u>datablock leak</u>	X	✓	X	
IC	<u>garbage data.</u>	✓	✓	X	!!
IC	<u>data bitmap / block contention</u>	✓	X	✓	*
	<u>leak / loose datablocks</u>	X	✓	✓	

① fsck : file system consistency check.

- reactive process.
- given a device & the file system on it, checks for consistency violations by reading the disk.
- superblock : FS size  $\leq$  device size  
checksum integrity check on the SB  
consistency of replicas.
- data blocks : are all datablocks accounted for  
leaks
  - some inode pts to the block.
  - if no pts. block bitmap can be flipped.
- inode state : all inodes in use have inode. bitmap set.
- links : ref cnt. in inode matches  
# files/entries that point to inode.

- duplicates: blocks shared in inodes.  
(need check for explicit setup or ~~viol~~ violation).

- bad blocks: out of range constants.

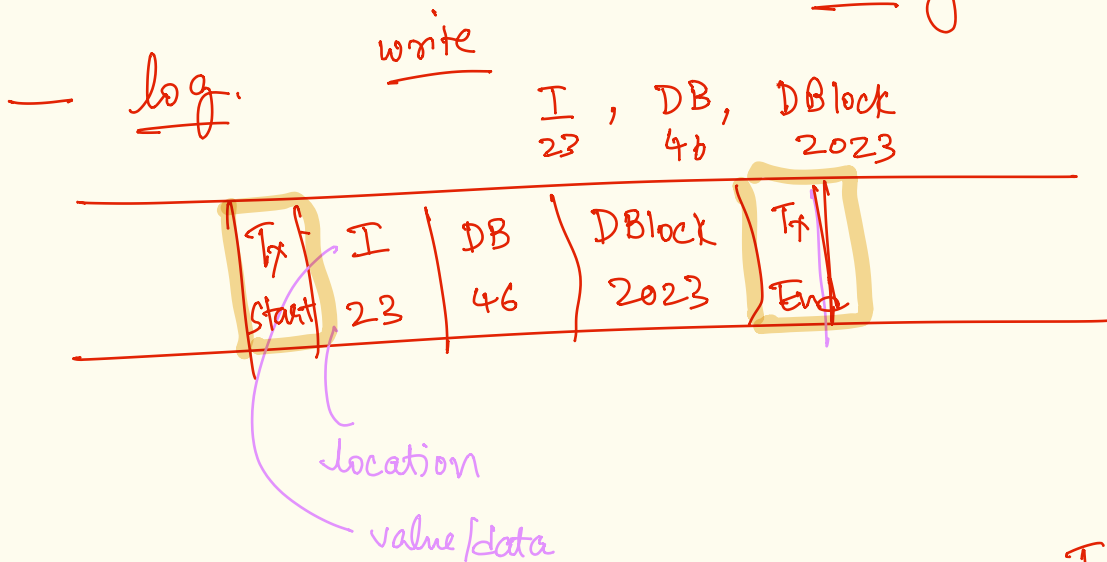
- directory checks — loops.

— ∴ are first 2 entries.

## ② write-ahead logging / journaling.

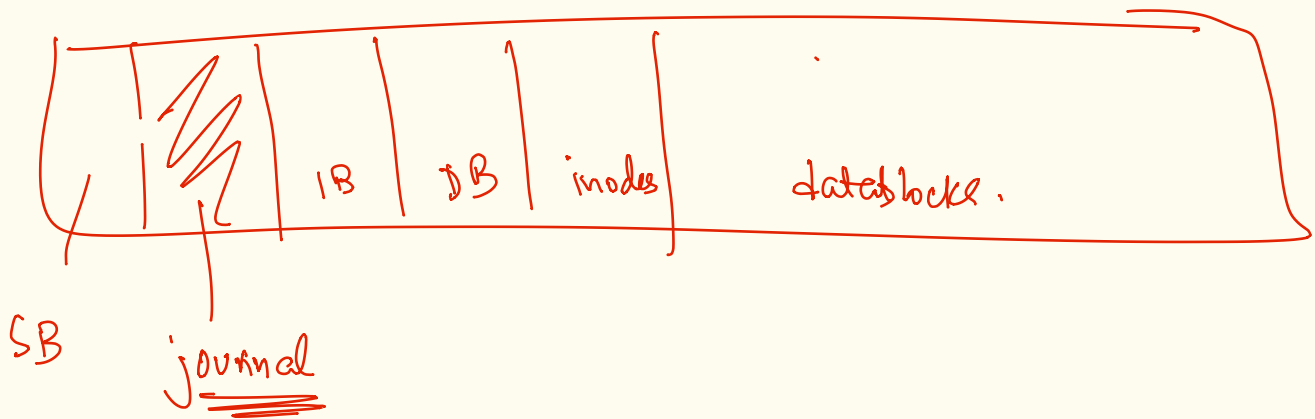
- information ~~is~~ about updates is written / logged separately.

- the log has transaction semantics.  
↳ all or nothing



— on (every) write . Tx updates Tx  
start end

- (i) update / write log to the journal
- (ii) commit ————— end of transaction.
- (iii) Check point — flush journal updates to disk (on actual locations)
- (iii) free / reuse the log.



③ RAID - Redundant Array of Inexpensive independent Disks.